

Correlation and Substitution in SPARQL

Daniel Hernández Claudio Gutierrez Renzo Angles

Version July 13, 2016

Abstract

In the current SPARQL specification the notion of correlation and substitution are not well defined. This problem triggers several ambiguities in the semantics. In fact, implementations as Fuseki, Blazegraph, Virtuoso and rdf4j assume different semantics.

In this technical report, we provide a semantics of correlation and substitution following the classic philosophy of substitution and correlation in logic, programming languages and SQL. We think this proposal not only gives a solution to the current ambiguities and problems, but helps to set a safe formal base to further extensions of the language.

This work is part of an ongoing work of Daniel Hernandez. These anomalies in the W3C Specification of SPARQL 1.1 were detected early and reported no later than 2014, when two erratas were registered (cf. <https://www.w3.org/2013/sparql-errata#errata-query-8> and <https://www.w3.org/2013/sparql-errata#errata-query-10>).

1 Introduction

The first version of this technical report served as a starting point to restart the discussion about the substitution and correlation in SPARQL. This issue was discussed in several threads on the W3C public-sparql-dev mailing list (see messages of Jun, 2016 in the mailing list archives¹) and a W3C Community Group² was created to discuss and address problems with the specification of the EXISTS clause in SPARQL.

This new version of the report fixes errors of the previous one and includes a formalization of two alternative semantics that are currently implemented: The first by Blazegraph and Fuseki, a semantics where substitution is never applied because variables that are not projected to resulting solutions are not visible from outside. The second, by Virtuoso and rdf4j, where every variable that is not projected to resulting solution are visible from outside, so they can be substituted. Within the same formal framework, we show the semantics presented in the previous technical report, where some variables are visible and other are not.

The main idea of the formal framework work as follows. Given a graph pattern P and a solution mapping μ , the Standard Spec. of SPARQL introduces the notion of $\text{substitute}(P, \mu)$, that is used to evaluate nested patterns. However, as we show in this report, this function substitute is not well defined and is contradictory with other parts of the specification. In this tech report, we define a similar function, which we call bind (to avoid clash names), that solves the problems found. It basically normalizes the pattern P before applying the mapping μ , giving a new structure $\text{norm}(P)$ that essentially renames variables so that each one plays the same role in every occurrence.

Structure of this technical report Section 2 presents an example of correlation using substitution to exemplify ambiguities of the current specification and differences of implementations. Section 3 describes the problem with current notion of substitution. In Section 4 we propose three alternative ways to define $\text{bind}(P, \mu)$ based in alternative definitions for $\text{norm}(P)$. Section 5 discuss how the proposed semantics are safe regarding with the use of blank nodes. Finally, in Section 6 we present several examples that illustrated how correlation is evaluated in each semantics and how implementations match them.

2 Evaluation of correlated variables

Consider the following simple SPARQL query that selects people of country j that have children, and consider as data the RDF graph depicted in Figure 1 below.

Listing 1

¹ <https://lists.w3.org/Archives/Public/public-sparql-dev/2016AprJun/>

² <https://www.w3.org/community/sparql-exists/>

```

1 SELECT ?parent
2 WHERE { ?parent :country :j
3         FILTER ( EXISTS { SELECT ?child
4                           WHERE { ?child :parent ?parent } }) }

```

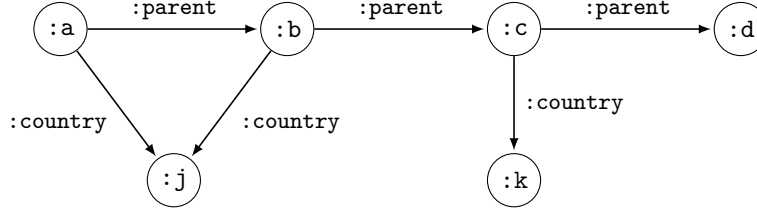


Fig. 1: RDF graph.

The engines Fuseki and Blazegraph give as solution two mappings: μ_a and μ_b , where μ_a is $\{?parent \mapsto :a\}$ and μ_b is $\{?parent \mapsto :b\}$. Virtuoso and rdf4j (formerly Sesame), on the other hand, give a different result, only μ_b . Why these differences?

What does the W3C Spec. tell? The query has the form `SELECT ?parent WHERE {P FILTER (EXISTS {Q})}`, where we will call P and Q respectively the outer and the inner graph patterns. How to evaluate this query? The W3C Spec. gives these two definitions that are relevant for this:

Definition 1 (Standard substitution, W3C Spec., §18.6). Let μ be a solution mapping and P be a graph pattern. Then, $\text{substitute}(P, \mu)$ is the graph pattern formed by replacing, for each x in $\text{dom}(\mu)$, every occurrence of a variable x in P by $\mu(x)$.

Definition 2 (Evaluation of Exists, W3C Spec. §18.6). Let μ be the current solution mapping for a filter and P a graph pattern: The value $\text{exists}(P)$, given $D(G)$, is true if and only if $\text{eval}(D(G), \text{substitute}(P, \mu))$ is a non-empty sequence.

In definition 2 above, the argument $D(G)$ denotes that the evaluation is done against the dataset D using the graph G . For the sake of the readability, in that follows we use the notation $\llbracket P \rrbracket_D$ instead of $\text{eval}(D(G), P)$.

The result of the `P FILTER (EXISTS {Q})` clause, according to the W3C Spec., should be the set Ω such that:³

$$\Omega = \{\mu \in \llbracket P \rrbracket_D \mid \llbracket \text{substitute}(Q, \mu) \rrbracket_D \text{ is not empty}\}.$$

So, it seems that Virtuoso and rdf4j follow the standard here: First, evaluate the pattern P, which give mappings μ_a and μ_b , and for each of them, perform

³ We will avoid the multiplicities in this report because the problems we report are independent of having set or multiset semantics.

the replacement in Q . As $\llbracket \text{substitute}(Q, \mu_a) \rrbracket_D = \emptyset$ and $\llbracket \text{substitute}(Q, \mu_b) \rrbracket \neq \emptyset$, the final solution is the mapping μ_b .

In defense of Fuseki and Blazegraph, let us say that the W3C Spec. says in other place (§12. Subqueries): “Note that only variables projected out of the subquery will be visible, or in scope, to the outer query.” That is, the variable `?parent` inside the `WHERE` clause is not visible from outside, and thus, Q cannot be changed by any mapping μ (in the sense of Defn. 1). Thus the `FILTER (EXISTS {Q})` is true, and thus the two mappings μ_a and μ_b qualify as final solutions.

The situation becomes even more involved when one considers another way of writing the previous query. Consider now the query in listing 2.

Listing 2

```

1 SELECT ?parent
2 WHERE { ?parent :country :j
3         FILTER ( EXISTS { SELECT ?child
4                           WHERE { ?child :parent ?chparent
5                                   FILTER (?chparent = ?parent) } }) }
```

The listings 3 and 4 present the queries resulting of applying the substitution on the inner graph pattern with the mappings μ_a and μ_b , respectively.

Listing 3

```

1 SELECT ?child
2 WHERE {
3   ?child :parent ?chparent
4   FILTER (?chparent = :a)
5 }
```

Listing 4

```

SELECT ?child
WHERE {
  ?child :parent ?chparent
  FILTER (?chparent = :b)
}
```

Only the second listing has solutions, so the evaluation of whole query returns $\{\mu_b\}$. Actually, in Virtuoso this query outputs the expected result, that is, $\{\mu_b\}$. On the contrary, in Fuseki and Blazegraph this query outputs no solutions. These systems are following another part of the W3C Spec. (12. Subqueries): “Due to the bottom-up nature of SPARQL query evaluation, the subqueries are evaluated logically first, and the results are projected up to the outer query.” Thus, they probably consider that the inner query returns error because there is a non-bound variable `?parent`.

Engines Examples presented in this report were tested in Fuseki 2.4.0⁴, Blazegraph Community Edition 2.1.0⁵, Virtuoso Open Source Edition 7⁶ and rdf4j 2.0 Milestone Builds⁷.

⁴ <https://jena.apache.org>

⁵ <https://www.blazegraph.com>

⁶ <https://github.com/openlink/virtuoso-opensource/tree/stable/7>

⁷ <http://rdf4j.org>

3 Problems with the current notion of substitution

The previous examples show that problems arise when it is not clear if occurrences of a variable are correlated⁸ and, in particular, if substitution has to be applied in a variable. In this section we will show that one of the main problems with nested queries in SPARQL is that the notion of *substitution* is not well defined. We will present a solution to this issue, which in turn helps to fix the whole semantics of nesting.

First, consider the simple graph pattern Q :

```
SELECT ?x WHERE { :a :p ?x }
```

and let μ be the solution mapping $\{?x \mapsto 1\}$. Then, $\text{substitute}(Q, \mu)$, understood literally from the standard, means replacing every occurrence of $?x$ with 1, that gives:

```
SELECT 1 WHERE { :a :p 1 }
```

Thus, the substitution method of the W3C Spec. is not well defined because it breaks the grammar of the **SELECT** clause.⁹

The notion of substitution was already present in SPARQL 1.0 in patterns of the form **P FILTER (C)**. In SPARQL 1.0 **C** is a Boolean clause, and here the substitution works fine because every occurrence of a variable could be replaced without breaking the grammar. The only case that required a special treatment was the function `bound(?x)` where $?x$ was not substituted, but checked if it was in the domain of the current solution.

On the contrary, in SPARQL 1.1, the clauses **EXISTS {Q}** and **NOT EXISTS {Q}** are filter constraints, thus allow nesting a graph pattern Q instead of the Boolean clause in a filter. And a graph pattern may contain variables with occurrences that are not replaceable (as we previously discussed) and variables with occurrences that are not “visible from outside” Q . Thus, the naive substitution, consisting on just replacing all occurrences of a variable, cannot be directly applied in SPARQL 1.1 as was in SPARQL 1.0.

4 Semantics of nested expressions with correlated variables

In this section we propose three alternatives for the function `bind`, that is defined as alternative to the function `substitute` for evaluating nested queries with correlated variables, two of which represent approaches existing in current implementations and the other the approach that was proposed in the previous version of this report.

⁸ We use “correlated variables” and “correlation” to indicate the occurrence of a variable x in and expression E whose value depends on the value of the occurrence of same variable x in an expression containing E . The paradigmatic occurrence of correlation in SPARQL is the expression **Q EXIST FILTER (P)**.

⁹ This and more subtle problems that a naive notion of substitution brings are well known in logic and algebra long ago. For example, a variable x cannot be substituted by a constant in all its occurrences in the first order formula $\forall x p(x)$ or in an expression like $\sum_{x \in A} (x + a)$.

4.1 The domain of a graph pattern

Given a graph pattern P , we denote as $\text{var}(P)$ to the set of variables that occur in P .

An interesting subset of $\text{var}(P)$ is the one that includes the variables that occur in the solutions of P , that we call the domain of P and denote as $\text{dom}(P)$, that is formally defined as follows¹⁰:

$$\text{dom}(P) = \{x \in \text{dom}(\mu) \mid \text{exists dataset } D \text{ with } \mu \in \llbracket P \rrbracket_D.\}$$

This definition of dataset cannot be used directly to compute the domain of a graph pattern, because requires the verification in all possible datasets. The following lemma shows that it is possible to give a method to compute the domain of a graph pattern using only its syntax.

Lemma 1 (In-domain variables). *Given a graph pattern P and a variable $?x$ occurring in P , then $?x \in \text{dom}(P)$ if and only if:*

1. *P is a basic graph pattern and $?x$ occurs in P .*
2. *If P is $Q \circ R$ where \circ is ‘.’, UNION or OPTIONAL and $?x \in \text{dom}(Q) \cup \text{dom}(R)$.*
3. *If P is Q MINUS R where \circ is ‘.’, and $?x \in \text{dom}(Q)$.*
4. *If P is GRAPH $?x$ $\{Q\}$.*
5. *If P is GRAPH u $\{Q\}$ and $?x \in \text{dom}(Q)$.*
6. *If P is VALUES (X) $\{B\}$ and $?x$ occurs in the list of variables X .*
7. *If P is Q BIND $(E$ AS $?x)$.*
8. *If P is Q BIND $(E$ AS $?y)$ and $?x \in \text{dom}(Q)$.*
9. *If P is Q FILTER (C) and $?x \in \text{dom}(Q)$.*
10. *If P is SERVICE u $\{Q\}$ and $?x \in \text{dom}(Q)$.*
11. *If P is SELECT X WHERE $\{Q\}$ if X is a list such that one of its elements is $?x$ or has the form $(E$ AS $?x)$.*

Note that in the SPARQL specification variables that are in the domain of graph pattern are called *in-scope* and also defined using the syntax (see 18.2.1 in the W3C Spec.). In this report we call them *in-domain* to stress the idea that they define the domain of the output.

¹⁰ Recall that $\text{dom}(\mu)$ is the domain of variables of μ (those where μ is defined) when μ is considered as a partial function over the set of all variables of the universe.

4.2 Syntax and variables roles

A variable in the syntax can play several roles. Two relevant ones are the role of representing the the output of a computation (output role) and the one representing the output of a previous computation (input role). For example, in the expression `let x be f(y) in { g(x,y) }` in a functional language, the variable `x` is playing the output role in the outermost occurrence and the input role in the innermost occurrence. On the other hand, `y` plays the input role in both occurrences.

These roles are crucial to understand how the substitution of variables by values work. Indeed, when a variable is in the input role, we can substitute it without breaking the syntax of the language. On the contrary, a variable in the output role cannot be substituted, because values cannot be used to name results of computations.

The question that arises is if we can distinguish the role of a variable occurrence in SPARQL. To answer this question, let us to consider the following types of syntactic constructs in which variables occur:

Expression. In a comparison (e.g., `?x < 2`), an scalar operation (e.g., `1+x`) or an scalar function (e.g., `substr(?x, 4)`).

Pattern. In a basic graph pattern (e.g., `?x :p ?y`).

Naming. In any place that only variables are allowed (e.g., `E AS ?x`), except when they occur in the `bound(.)` function (e.g., `bound(?x)`), that is an special case.

The occurrence of variables in the three types of constructs are associated with the output or input roles as shown in Table 1.

	Input	Output
Expression	×	
Pattern	×	×
Naming		×

Tab. 1: Possible variable roles in types of syntactic constructs.

In occurrences in expressions, it is clear that the variable can be substituted, because it refers to a value to be used inside the expression. Similarly, in naming occurrences, it is clear that the variable cannot be substituted, because breaks the grammar. Moreover, the variable will be used to refer the value in a future computation so its name cannot be forgotten nor changed.

In the case of occurrences in patterns the variable could have both roles. Indeed, we can substitute the variable with a value without breaking the semantics, so the variable is playing the input role. On the other hand, if the variable is not substituted, then it will bind a value from the data that will be available for future computations, so it is playing an output role.

The substitution of variables that are in syntactic constructs of type pattern (i.e., in a basic graph pattern) has another issue: A variable $?x$ that is replaced by a value in a basic graph pattern P does not appear in the solutions of evaluating P . Thus, the domain of P will be reduced after the substitution.

This, reduction in the domain of a graph pattern after a substitution, may produce odd results. Indeed, let P and Q be respectively the basic graph patterns $\{?x : p ?y\}$ and $\{?y : p ?z\}$. Let P' and Q' be the results of substituting $?y$ by $:b$ in P and Q , respectively. Let μ be the solution $\{?y \mapsto :b\}$. Then, the graph pattern $P.Q$ has less solutions than $P'.Q'$ over the dataset $\{(:a, :p, :b), (:b, :p, :c)\}$. This contradicts, the intuition that substituting variables with values restrict the results.

4.3 Normalization

The normalization of a graph pattern or expression P is defined to avoid variables with role ambiguity (i.e., that has simultaneously input and output roles) by changing the structure of P and replacing every variable occurring in P with a different fresh variable for each scope that can be determined for the variable. After the normalization process, variables that can be substituted will occur only in syntactical constructs of type expression, so solving the issue described at the end of the previous section.

Definition 3 (Normalization). The normalization of the pattern P , that we denote as $\text{norm}(P)$, is a triple (P', d, g) , where P' is a pattern whose variables must be all fresh and d and g are partial functions whose domain and ranges are as follows:

$$\begin{aligned} d &: \text{var}(P') \rightarrow \text{dom}(P), \\ g &: \text{var}(P') \rightarrow \text{var}(P), \end{aligned}$$

d is surjective and the domains of d and g are disjoint.

Intuitively, d and g are functions that associate (record) the correspondence of the fresh variables of P' with the corresponding original variables P . The function d represents occurrences of variables that are in the solutions of P and g represents occurrences of variables that can be substituted by values that μ maps. The sets $\text{range}(d)$ and $\text{range}(g)$ could have elements in common. For example, if P is the graph pattern $Q.R$ then a variable can be in the domain of Q and simultaneously be a global variable in R ,

To give an intuition, here there is an illustration of a normalization in a simple case. Let P be:

$$\{ :a :p ?x \} . \{ :b :q ?y \text{ FILTER } (?y < ?x) \}.$$

The result of normalizing P is (P', d, g) where P' is

$$\{ :a :p x_1 \} . \{ :b :q y_1 \text{ FILTER } (y_1 < x_2) \}$$

and d and g are respectively the functions

$$\begin{aligned} d &:= \{x_1 \mapsto ?x, y_1 \mapsto ?y\}, \\ g &:= \{x_2 \mapsto ?x\}. \end{aligned}$$

(We use a different notation for variables P' to stress the idea that they are fresh.) Note that in the pattern P' in (P', d, g) each variable plays a unique role, and the functions d, g “tell” what is the role of each variable and their relationships.

Note that this example uses a particular normalization according with a specific semantics. An alternative semantics may produce a different normalization (which is only designed to make the role of each variable independent of its occurrence).

4.4 Substitution and correlated evaluation

We need a pair of notations before introducing the main notions. Given a partial function f that maps variables to variables and an structure of expression A where some of this variables occur, then $f(A)$ denotes the result of renaming consistently in A every variable $x \in \text{dom}(f)$ by $f(x)$. Functions can be viewed a set of ordered pairs. We will use the notation $x \mapsto f(x)$ instead of $(x, f(x))$ to stress the notion of mapping. Thus, the symbol \emptyset (used commonly to denote empty sets) also denotes empty functions.

Now we are ready to present our main notion:

Definition 4 (Mapping substitution). Let P be a graph pattern, μ a solution mapping and d and g be functions that map variables to variables. Then $\mu(P, d, g)$ is the graph pattern $d(P')$, where P' is the graph pattern resulting of the following substitutions in P :

1. For each binding $x \mapsto ?x$ in g substitute every occurrence of $\text{bound}(x)$ by **TRUE** if $?x \in \text{dom}(\mu)$ or by **FALSE** if $?x \notin \text{dom}(\mu)$.
2. Then, for each binding $x \mapsto ?x$ in d substitute every occurrence of x by $\mu(?x)$ if $?x \in \text{dom}(\mu)$ or by $?x$ if $?x \notin \text{dom}(\mu)$.

Definition 5 (Main: Correlated graph pattern or expression). Let P be a graph pattern or expression, μ be a solution and norm be a function that receives a graph pattern and returns triple (P', d, g) where P' is a graph pattern and d and g are functions that map variables to variables. Then:

$$\text{bind}(P, \mu) = \begin{cases} \mu(\text{norm}(P)) \cdot \mu|_{\text{dom}(P)} & \text{if } P \text{ if a graph pattern,} \\ \mu(\text{norm}(P)) & \text{if } P \text{ if an expression.} \end{cases}$$

Note that $\mu|_{\text{dom}(P)}$ denotes the inline data that codify exactly the multiset containing the solution $\mu|_{\text{dom}(P)}$ with multiplicity 1. For example, if $\mu|_{\text{dom}(P)}$ is the solution $\{?x \mapsto 1, ?y \mapsto 2\}$ then it is codified as **VALUES** $(?x \ ?y) \{(1 \ 2)\}$.

The function `bind` can be used in any place where the function substitute is used by the Standard Spec. For example, given a dataset D , the graph patterns P and Q and the expression E , then:

$$\begin{aligned} \llbracket P \text{ FILTER } (\text{EXISTS } \{Q\}) \rrbracket_D &= \{\mu \in \llbracket P \rrbracket_D \mid \llbracket \text{bind}(Q, \mu) \rrbracket_D \text{ is not empty}\} \\ \llbracket P \text{ BIND } (E \text{ AS } ?x) \rrbracket_D &= \llbracket P \rrbracket_D \bowtie \{\{?x \mapsto \llbracket \text{bind}(E, \mu) \rrbracket_D\}\} \end{aligned}$$

In what follows, we present three variants of the normalization function. Each one, according to Definition 5 will give a particular semantics. Given a graph pattern P and its normalization (P', d, g) , these variants differ essentially in the variables occurring in P that are included in the range of the function g . Intuitively, variables that are excluded of the ranges of d and g can be considered local, because the normalization renames them to fresh variables and does not record the original names.

4.5 Semantics S1

According S1 all variables that are not in the domain of a graph pattern are considered local. Thus, the normalization in S1 is defined as follows:

Definition 6 (Normalization in S1). Given a graph pattern P , then $\text{norm}(P)$ is (P', d, \emptyset) where:

1. d is a surjective function that maps fresh variables to variables in $\text{dom}(P)$.
2. P' is $h(d^{-1}(P))$ where h is a function that maps variables in $\text{var}(P) \setminus \text{dom}(P)$ to fresh variables.

At the end of this procedure is ensure that all local variables in P are substituted P' with fresh variables that will be not substituted again because the third component of the normalization is empty. This is summarized in the following result.

Lemma 2. *According the semantics S1, given a graph pattern P , a solution mapping μ and a dataset D , then:*

$$\llbracket \text{bind}(P, \mu) \rrbracket_D = \llbracket P \rrbracket_D \bowtie \{\mu|_{\text{dom}(P)}\}$$

4.6 Semantics S2

Before defining this semantics we need some definitions that will help us in the notation.

Definition 7 (The filter clause). Given two variables $?x$ and $?y$ then $F_{?x?y}$ is the operator `FILTER` $(!(\text{bound}(?x) \ \&\& \ \text{bound}(?y)) \ || \ ?x = ?y)$.

The operator $F_{?x?y}$ help us to rewrite a variable that is in the domain of a graph pattern as a variable whose visibility is global according S2 and S3. For example, let P and Q be the graph patterns $\{a : p \ ?x\}$ and $\{x : p \ ?y\}$ $F_{?x?y}$, respectively. Then, intuitively $\llbracket P \rrbracket_\mu = \llbracket Q \rrbracket_\mu$ for every mapping μ .

Definition 8 (Consequently renaming). Let f and g be two functions that map variables to variables where g is injective. Then, $\text{cr}(f, g)$ is the function $g^{-1}|_A \cdot f|_A$ where A is $\text{range}(f) \cap \text{range}(g)$ and “ \cdot ” denotes the composition of functions¹¹.

If a graph pattern P is composed of the graph patterns Q and R , then results natural defining the normalization of P as a composition of the respective normalizations Q' and R' of its components. Because the normalization of these components are performed independently, the variables in the domain may be different, thus it is needed to rename variables in one of the components to make both renaming consequent in the outputs of both components. The following lemma show that given to renamings f and g where g is injective, then $\text{cr}(f, g)$ can be used to generate a function g' that is compatible with f , that is $(\text{cr}(f, g))(g)$.

Lemma 3. *Given two functions f and g that map variables to variables where g is injective, then $(\text{cr}(f, g))(g|_A) = f|_A$.*

At this point we are ready to proceed with the formalization of the normalization in the semantics S2.

Definition 9 (Normalization in S2). Given a graph pattern or expression P , then:

1. If P is a basic graph pattern, then $\text{norm}(P)$ is $(d^{-1}(P), d, \emptyset)$, where d is a function that contains a binding $x \mapsto ?x$ for every variable $?x$ in $\text{var}(P)$.
2. If P is **SELECT X WHERE {Q}** (where X is a list of variables), then $\text{norm}(P)$ is (P', d_P, g_P) , where:

$$\begin{aligned} P' &= \text{SELECT } X' \text{ WHERE } \{Q'\} \\ (Q', d_Q, g_Q) &= \text{norm}(Q) \\ d_P &= d_Q|_{d_Q^{-1}(\text{dom}(P))} \\ g_P &= g_Q \\ X' &= d_P^{-1}(X) \end{aligned}$$

Note that $d_Q^{-1}(\text{dom}(P))$ is the preimage in d_Q of $\text{dom}(P)$. That is, the set of variables used in Q to rename variables that are projected in the solution of P . Thus, $d_Q|_{d_Q^{-1}(\text{dom}(P))}$ is the renaming of variables used in Q , restricted to the domain of P . Similarly, $d_P^{-1}(X)$ is renaming of variables in X that is consequent with the renaming done in Q .

Bindings $x \mapsto ?x$ in d_Q such that $?x$ is not in the domain of P are not included in d_P nor in g_P . This, is interpreted as that the cocurrences of $?x$ associated to this bindings are assumed local.

¹¹ Note that that is if $x \in \text{dom}(f)$ and $f(x) \in \text{dom}(g)$ then $(f \cdot g)(x) = g(f(x))$.

3. If P is $Q \circ R$ where \circ is \cdot , **OPTIONAL** or **UNION**, then $\text{norm}(P)$ is (P', d_P, g_P) , where:

$$\begin{aligned} P' &= Q' \circ f(R') \\ (Q', d_Q, g_Q) &= \text{norm}(Q) \\ (R', d_R, g_R) &= \text{norm}(R) \\ d_P &= d_Q \cup f(d_R) \\ g_P &= g_Q \cup g_R \\ f &= \text{cr}(d_Q, d_R) \end{aligned}$$

Note that f is a renaming that ensure that the normalizations of Q and R use the same common domain variables when they are combined.

4. If P is $Q \text{ MINUS } R$, then $\text{norm}(P)$ is $(Q' \text{ MINUS } f(R'), d_P, g_P)$, where:

$$\begin{aligned} (Q', d_Q, g_Q) &= \text{norm}(Q) \\ (R', d_R, g_R) &= \text{norm}(R) \\ d_P &= d_Q \\ g_P &= g_Q \cup g_R \\ f &= \text{cr}(d_Q, d_R). \end{aligned}$$

The function f is a renaming that ensures that the variables used in the variables that are common in domains of the normalizations of Q and R are renamed to the the same fresh variables.

Note that variables that are in $\text{dom}(R) \setminus \text{dom}(Q)$ are replaced with fresh variables that are not included in the domains of d_P and g_P . Thus, they are assumed local.

5. If P is **GRAPH** $u \{Q\}$ where u is an IRI, then $\text{norm}(P)$ is $(\text{GRAPH } u \{Q'\}, d, g)$, where $\text{norm}(Q) = (Q', d, g)$.
6. If P is **GRAPH** $?x \{Q\}$, then $\text{norm}(P)$ is $(\text{GRAPH } x \{Q'\}, d_P, g_P)$, where:

$$\begin{aligned} (Q', d_Q, g_Q) &= \text{norm}(Q) \\ d_P &= \begin{cases} d_Q & \text{if } ?x \in \text{range}(d) \\ d_Q \cup \{x \mapsto ?x\} & \text{otherwise } (x \text{ is fresh}) \end{cases} \\ g_P &= g_Q. \end{aligned}$$

7. If P is **SERVICE** $u \{Q\}$ where u is an IRI, then $\text{norm}(P)$ is (P', d, g) , where $\text{norm}(Q) = (Q', d, g)$ and $P' = \text{SERVICE } u \{Q'\}$.

8. If P is $Q \text{ FILTER } (C)$, then $\text{norm}(P)$ is (P', d_P, g_P) where:

$$\begin{aligned} P' &= Q' \text{ FILTER } (f(C')), \\ (Q', d_Q, g_Q) &= \text{norm}(Q), \\ (C', \emptyset, g_C) &= \text{norm}(C), \\ d_P &= d_Q, \\ g_P &= g_Q \cup f(g_C), \\ f &= \text{cr}(d_Q, d_C). \end{aligned}$$

9. If P is $\text{VALUES } (X) \{B\}$ where X is a list of variables and B is a list of bindings to the variables, then $\text{norm}(P)$ is (P', d, \emptyset) , where d has a binding $x \mapsto ?x$ for each variable $?x$ in X and $P' = d^{-1}(P)$.

10. If P is $Q \text{ BIND } (E \text{ AS } ?x)$ then $\text{norm}(P)$ is (P', d_P, g_P) where:

$$\begin{aligned} P' &= Q' \text{ BIND } (f(E') \text{ AS } x) \\ (Q', d_Q, g_Q) &= \text{norm}(Q) \\ (E', \emptyset, g_E) &= \text{norm}(E) \\ d_P &= d_Q \cup \{x \mapsto ?x\}, \\ g_P &= g_Q \cup f(g_E), \\ f &= \text{cr}(d_Q, d_E). \end{aligned}$$

11. If P is an expression and $\{Q_1, \dots, Q_n\}$ is the set of graph patterns that are directly contained into maximal occurrences of **EXISTS** clauses in P (we say that an **EXISTS** clause occurrence i is maximal in P if does not occur another **EXISTS** clause j containing i in P). Then $\text{norm}(P)$ is (P_n, \emptyset, g_n) computed recursively as follows:

- (a) Let P_0 be P and g_0 be the function that include a binding $x \mapsto ?x$ for every variable $?x$ in P that does not occur in any of the graph patterns $\{Q_1, \dots, Q_n\}$.
- (b) For each Q_k of the graph patterns in the maximal **EXISTS** clauses, let (Q'_k, d'_k, g'_k) be $\text{norm}(Q_k)$. Then, let P_k the result of replacing in P_{k-1} the occurrence of Q_k by $f(Q_k) F_{x_1 y_1} \dots F_{x_m y_m}$ where:

$$\begin{aligned} f &= \text{cr}(g_0, g'_k), \\ g_k &= g_{k-1} \cup f(g'_k) \cup h(d'_k) \end{aligned}$$

and h be a function $\{y_1 \mapsto x_1, \dots, y_m \mapsto x_m\}$ that has a binding $y \mapsto x$ for each binding $x \mapsto ?x$ in d'_k .

4.7 Semantics S3

In the rules 2 and 4, the semantics S2 assumes that pattern occurrences that are not in the domain of a graph pattern are local, so they are not included in the

global bindings. On the contrary, S3 assumes that they are global, so the query is modified to move variables from graph occurrences to expression occurrences using operations F_{xy} . After this transformation, substitution can be applied in the same way that in semantics S2.

Definition 10 (Normalization in S3). Given a graph pattern or expression P, then the normalization of P is computed using the same rules enumerated in the definition of the normalization of S3, except the rules 2 and 4, that are replaced with the following rules:

2. If P is **SELECT X WHERE {Q}** (where X is a list of variables), then $\text{norm}(P)$ is (P', d_P, g_P) , where:

$$\begin{aligned} P' &= \text{SELECT } X' \text{ WHERE } \{Q' F_{x_1 y_1} \dots F_{x_m y_m}\} \\ (Q', d_Q, g_Q) &= \text{norm}(Q) \\ d_P &= d_Q|_{\text{dom}(P)} \\ g_P &= g_Q \cup h(d_Q|_{\text{dom}(Q) \setminus \text{dom}(P)}) \\ X' &= d_P^{-1}(X) \end{aligned}$$

and h is a function $\{y_1 \mapsto x_1, \dots, y_m \mapsto x_m\}$ that has a binding $y \mapsto x$ for each binding $x \mapsto ?x$ in $d_Q|_{\text{dom}(Q) \setminus \text{dom}(P)}$.

4. If P is **Q MINUS R**, then $\text{norm}(P)$ is (P', d_P, g_P) where:

$$\begin{aligned} P' &= Q' \text{ MINUS } \{f(R') F_{x_1 y_1} \dots F_{x_m y_m}\} \\ (Q', d_Q, g_Q) &= \text{norm}(Q) \\ (R', d_R, g_R) &= \text{norm}(R) \\ d_P &= d_Q \\ g_P &= g_Q \cup g_R \cup h(d_R|_{\text{dom}(R) \setminus \text{dom}(P)}) \\ f &= \text{cr}(d_Q, d_R) \end{aligned}$$

and h is a function $\{y_1 \mapsto x_1, \dots, y_m \mapsto x_m\}$ that has a binding $y \mapsto x$ for each binding $x \mapsto ?x$ in $d_R|_{\text{dom}(R) \setminus \text{dom}(P)}$.

5 Substitution and blank nodes

Peter F. Patel-Schneider noticed in the W3C mailing list of SPARQL that substitution has problems with blank nodes and a semantics where every variable can be substituted. Let P be the inner graph pattern of a query and $?x$ be a variable that can be substituted in P, in particular with a blank node $_:b$. Then, there are the following options:

1. $?x$ occurs in a basic graph pattern and is substituted by $_:b$. Then, according the specification $_:b$ is interpreted as an existential variable and scoped to the basic graph pattern. Thus, $_:b$ may represent any element on the graph, not only $_:b$.

2. $?x$ is in an expression occurrence. Then, the substitution of $?x$ by $_:b$ restricts the resulting bindings to those where $?x$ is bound to $_:b$.
3. $?x$ is in the domain of P . Then, results of P that are not compatible with $\{?x \mapsto _:b\}$ are discarded.

The first case results contradictory with the other two. In fact, it does not restrict the variable $?x$ as the other do (and as it is expected for substitution).

In the semantics S3 proposed in this technical report a variable $?x$ occurring in a pattern occurrence is considered replaceable with values that come from the current solution mapping. However, before this substitution the normalization process moves $?x$ from the pattern occurrence to an expression occurrence using a renaming of $?x$ to $?y$ and then using an operator $F_{?x?y}$. Thus, we can conclude the following lemma:

Lemma 4. *The semantics S1, S2 and S3 are safe respect with the blank nodes substitution issue.*

6 Correlation in implementations

This section presents examples of how the proposed semantics work and how different implementations match them. Queries presented in this section are run against the RDF graph depicted in Figure 1.

For each query, the actual results given by each implementation is shown at the end of the section.

Example 1

```
SELECT ?parent
WHERE { ?parent :country :j
        FILTER ( EXISTS { ?child :parent ?parent } ) }
```

This query gets people of country :j having children. That is, select people that has solutions for the inner query. The variable **?people** is in-domain in the inner graph pattern for all semantics. Thus, the results of inner graph pattern are filtered to be compatible with solutions of the outer query μ_a and μ_b . The only solution that has results for the inner graph pattern is μ_b in each of the three semantics S1, S2 and S3.

Example 2

```
SELECT ?parent
WHERE { ?parent :country :j
        FILTER ( EXISTS { SELECT ?child
                           WHERE { ?child :parent ?parent } } ) }
```

This query is similar to the presented in Example 1. However, in this case the variable **?parent** in the inner query is in local according with S1 and S2, and


```

FILTER (?parent = 1 ||
        ?parent != 1 )}}})

```

In this example, the filter clause `?parent = 1 || ?parent != 1` is a tautology when `?parent` is bound. Otherwise, both sides of the disjunction are evaluated as error so the whole clause gets an error. Thus, the output of this query is $\{\mu_a, \mu_b\}$ according S2 and S3 and $\{\}$ according S1.

Example 7

```

SELECT ?parent
WHERE { ?parent :country :j
        FILTER ( EXISTS { SELECT *
                           WHERE { ?child :parent ?chparent
                                    FILTER (?parent = 1 ||
                                             ?parent != 1 )}}})

```

This example is equivalent with Example 6. In fact, the semantics of the wildcard ‘*’ is the list of all variables that are in-domain of the query.

Example 8

```

SELECT ?parent
WHERE { ?parent :country :j
        FILTER ( EXISTS { SELECT ?child
                           WHERE { ?child :parent ?parent
                                    FILTER (?parent = :c)}}})

```

In this query the variable `?parent` has three occurrences. The first is in the outer graph pattern and the other two in the inner graph pattern. In any of the semantics `?parent` is bound to `:a` and `:b` in the solutions of the outer graph pattern.

According S2 and S3 the variable `?parent` is local in the inner graph pattern. Thus it is bound to `:a`, `:b` and `:c` in the inner graph pattern. Then, the filter clause of the inner graph pattern is true for `:c`. Thus, the result of this query is $\{\mu_a, \mu_b\}$.

On the other hand, according S3 the variable `?parent` is global in the inner graph pattern. So, it is replaced with the values coming from the outer graph pattern. None of this values satisfy the condition `?parent = :c`. Thus, the result of this query is $\{\}$.

Example 9

```

SELECT ?parent
WHERE { ?parent :country :j
        FILTER ( EXISTS { SELECT ?child
                           WHERE { ?child :parent ?parent
                                    FILTER (EXISTS{?parent :parent :d}})})

```

On the dataset used in these examples, This query seems to be equivalent to the previous query (presented in Example 8), because the graph pattern `?parent :parent :d` only has solutions if `?parent` is `:c`. Thus, in S2 and S3 the result of this query is $\{\}$. On the other hand, in S1 the result of this query is $\{\mu_a, \mu_b\}$.

Example 10

```
SELECT *
WHERE { { { ?x :p ?y } OPTIONAL { ?y :q ?z } }
        FILTER ( EXISTS { ?z :r ?v } ) }
```

Consider the following RDF graph that is the dataset D where we will evaluate this query.

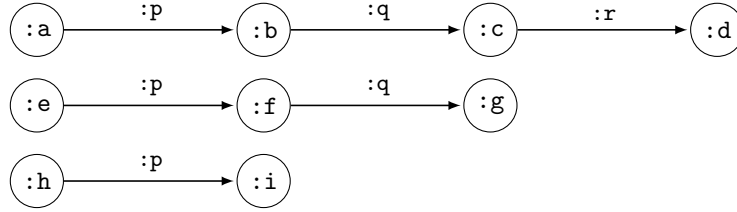


Fig. 2: RDF graph.

The normalization of the inner graph pattern gives the same result in the three semantics, because the variables `?z` and `?v` are in the domain of the inner graph pattern. Thus, the evaluation of this query is the set Ω defined as:

$$\Omega = \{\mu \in \llbracket P \rrbracket_D \mid \llbracket Q \rrbracket_D \bowtie \{\mu\} \text{ is not empty}\},$$

where P and Q are the outer and inner graph patterns, respectively. $\llbracket P \rrbracket_D$ and $\llbracket Q \rrbracket_D$ are respectively the sets $\{\mu_{abc}, \mu_{efg}, \mu_{hi}\}$ and $\{\mu_{cd}\}$ where:

$$\begin{aligned} \mu_{abc} &= \{?x \mapsto :a, ?y \mapsto :b, ?z \mapsto :c\}, \\ \mu_{efg} &= \{?x \mapsto :e, ?y \mapsto :f, ?z \mapsto :g\}, \\ \mu_{hi} &= \{?x \mapsto :h, ?y \mapsto :i\}, \\ \mu_{cd} &= \{?z \mapsto :c, ?v \mapsto :d\}. \end{aligned}$$

Thus, $\Omega = \{\mu_{abc}, \mu_{hi}\}$.

Summary

The following table summarizes the results that the example queries get for each of the studied semantics and the results that the studied implementations actually output.

#	S1	S2	S3	rdf4j	Virtuoso	Fuseki	Blazegraph
1	$\{\mu_b\}$	$\{\mu_b\}$	$\{\mu_b\}$	$\{\mu_b\}$	$\{\mu_b\}$	$\{\mu_b\}$	$\{\mu_b\}$
2	$\{\mu_a, \mu_b\}$	$\{\mu_a, \mu_b\}$	$\{\mu_b\}$	$\{\mu_b\}$	$\{\mu_b\}$	$\{\mu_a, \mu_b\}$	$\{\mu_a, \mu_b\}$
3	$\{\}$	$\{\mu_b\}$	$\{\mu_b\}$	$\{\mu_b\}$	$\{\mu_b\}$	$\{\}$	$\{\}$
4	$\{\}$	$\{\mu_a, \mu_b\}$	$\{\mu_a, \mu_b\}$	$\{\mu_a, \mu_b\}$	$\{\}$	$\{\}$	$\{\}$
5	$\{\}$	$\{\mu_b\}$	$\{\mu_b\}$	$\{\mu_b\}$	$\{\mu_b\}$	$\{\}$	$\{\}$
6	$\{\}$	$\{\mu_a, \mu_b\}$	$\{\mu_a, \mu_b\}$	$\{\mu_a, \mu_b\}$	$\{\mu_a, \mu_b\}$	$\{\}$	$\{\}$
7	$\{\}$	$\{\mu_a, \mu_b\}$	$\{\mu_a, \mu_b\}$	$\{\mu_a, \mu_b\}$	$\{\mu_a, \mu_b\}$	$\{\mu_a, \mu_b\}$	$\{\}$
8	$\{\mu_a, \mu_b\}$	$\{\mu_a, \mu_b\}$	$\{\}$	$\{\}$	$\{\mu_a, \mu_b\}$	$\{\mu_a, \mu_b\}$	$\{\mu_a, \mu_b\}$
9	$\{\mu_a, \mu_b\}$	$\{\mu_a, \mu_b\}$	$\{\}$	$\{\}$	$\{\}$	$\{\mu_a, \mu_b\}$	$\{\mu_a, \mu_b\}$
10	$\{\mu_{abc}, \mu_{hi}\}$	$\{\mu_{abc}, \mu_{hi}\}$	$\{\mu_{abc}, \mu_{hi}\}$	$\{\mu_{abc}, \mu_{hi}\}$	$\{\mu_{abc}\}$	$\{\mu_{abc}, \mu_{hi}\}$	$\{\mu_{abc}, \mu_{hi}\}$

We distinguish two groups of implementations. Blazegraph and Fuseki match S1 and Virtuoso and rdf4j match S3. However, only Blazegraph and rdf4j match their respective semantics in all the examples.

Fuseki agree with S1 except in the query of Example 7 with . The queries in the examples 6 and 7 are equivalent. However, in Fuseki the results differ. This seems as a bug of Fuseki.

Virtuoso agree with S3 except in the queries of examples 4, 8 and 10. The result given in the query 4 seems as a bug, because it is contradictory that `?parent` is unbound when it is bound in query 5. The result given in the query 8 require more study because it is not clear if it is a bug or it is motivated by a different interpretation of the correlation. The result given in the query 10 shows that there are an special treatment with unbound variables.